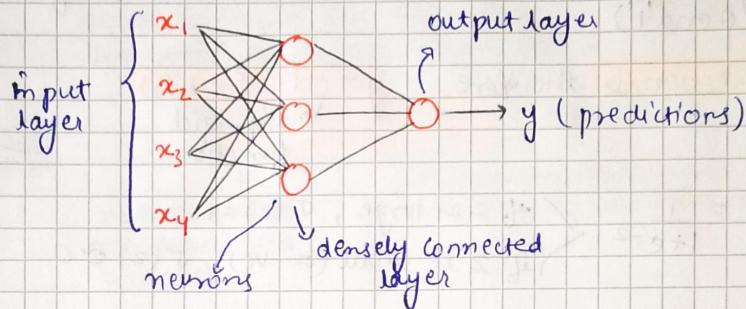


## Neural networks and deep learning

- ① **Neuron** → Within an artificial neural network, a "neuron" is a mathematical function that models the functioning of a biological neuron.  
 → Typically, a neuron compute the weighted average of its input, and this sum is passed through a non-linear function, often called "activation function".



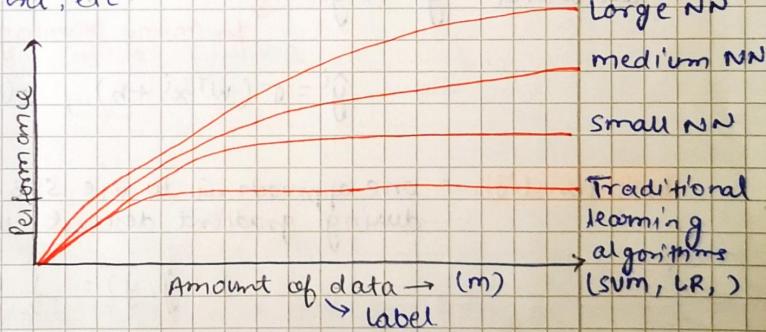
- ② **Supervised learning** → Also known as supervised machine learning, is a subcategory of machine learning and AI. It is defined by its use of labelled datasets to train algorithms to classify data or predict outcomes accurately.

Structured data → Structured data, in terms of deep learning means that every feature has a well-defined meaning.

Unstructured data → Unstructured data can be categorised in Audio signals, image, text, etc.

- ③ **Deep learning progress** →

→ Deep learning is progressing by scaling. Scaling meaning with the size of the network as well as size of the data



- ④ **Logistic Regression** → It can be defined as an algorithm used for binary classification

Binary classification → If we have an RGB image, we make an input vector having the dimensions  $(m \times n \times 3)$  where  $m =$  e.g.  $64 \times 64, 128 \times 128$ , etc.

**Notation** ① Single training example →  $(x, y) \rightarrow y \in \{0, 1\} \rightarrow \text{labels}$

② Training set  $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\} \rightarrow x \in \mathbb{R}^{n \times m} \}$  input vector

$m_{\text{train}} \rightarrow \text{no. of training examples}, m_{\text{test}} \rightarrow \text{no. of test examples}$

$$\text{③ } X = \begin{bmatrix} 1 & 1 & 1 \\ x^1 & x^2 & \dots & x^m \\ 1 & 1 & 1 \end{bmatrix} \quad \begin{matrix} \uparrow \\ m_x \end{matrix}$$

$X \in \mathbb{R}^{m_x \times m}$   
 → matrix with the dimension  $m_x \times m$

$$\textcircled{1} \quad \mathbf{y} = [y^1, y^2, \dots, y^m] \quad \{ \mathbf{y} \in \mathbb{R}^{1 \times m} \}$$

→ Given an input vector  $\mathbf{x}$  (image), we want to predict whether it is a cat or not or

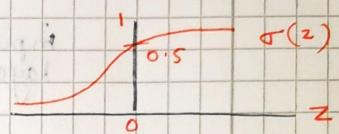
$$\hat{y} = P(y=1 | \mathbf{x}) \rightarrow \text{label "1" given the value of } \mathbf{x} \in \mathbb{R}^n$$

The parameters which we use are  $w \in \mathbb{R}^{n \times 1}$ ,  $b \in \mathbb{R}$

One idea is to use linear regression (i.e.  $\hat{y} = w^T \mathbf{x} + b$ ). But it seems it is not a very good approach since we want a chance of our output (i.e. a value between 0 and 1)

Logistic Regression is therefore

$$\hat{y} = \sigma(w^T \mathbf{x} + b) \quad \begin{matrix} z \\ \downarrow \text{Sigmoid function} \end{matrix}$$



$$\sigma(z) = \frac{1}{1+e^{-z}} \quad \begin{cases} \text{if } z \text{ is large, } \sigma(z) \approx 1 \\ \text{if } z \text{ is small (neg), } \sigma(z) \approx 0 \end{cases}$$

\* when implementing logistic regression, our goal should be to learn the parameters  $w$  and  $b$  so that  $\hat{y}$  becomes a good estimate of  $y$  being equal to 1

→ Logistic Regression cost function → To train the parameters "w" and "b" of our model, we define a cost function.

given  $\{\hat{y} = \sigma(w^T \mathbf{x} + b), \sigma(z) = \frac{1}{1+e^{-z}}\}$  and  $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$

we want  $\hat{y}^i \approx y^i$  for many training examples,

$$\hat{y}^i = \sigma(w^T x^i + b), \sigma(z^i) = \frac{1}{1+e^{-z^i}}$$

Loss function → one approach is to use SSE, but since it does not work good during gradient descent, we have to find another way.

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

A good approach is to use:  $L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log(1-\hat{y}))$

An approximation of this is that

① if  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y} \Rightarrow$  for  $y=1$ , we want  $-\log \hat{y}$  to be as ~~big~~<sup>small</sup> as possible meaning  $\hat{y}$  to be very big. But since we know that  $\hat{y}$  is a function of sigmoid, it can never be greater than 1.

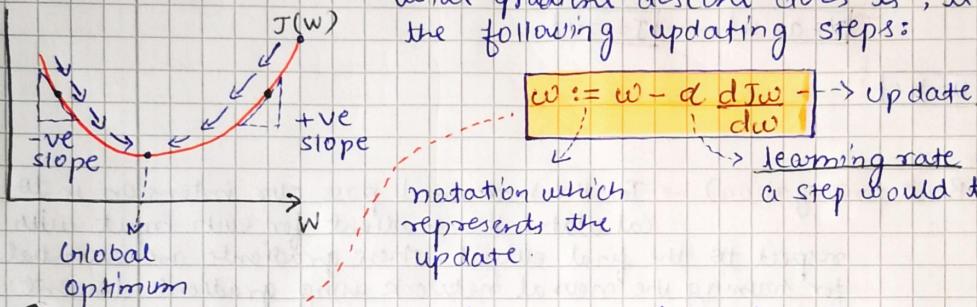
② if  $y=0$ :  $L(\hat{y}, y) = -\log(1-\hat{y}) \Rightarrow$  similarly, we want  $\log(1-\hat{y})$  to be large, and therefore we want  $\hat{y}$  to be as ~~small~~<sup>big</sup> as possible.

Cost function:  $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^i, y^i)$  over the entire training set.

→ **Gradient Descent** → It is an optimization algorithm which is used to minimize some function by iteratively moving in the direction of steepest descent as defined by the negative of the gradient. In machine learning, we use gradient descent to update the parameters ( $w, b$ ) of our model. Parameters refers to the "coefficients" in linear regression and weights in neural networks.

note \* In other words, we want to find  $w$  and  $b$  that minimizes  $J(w, b)$

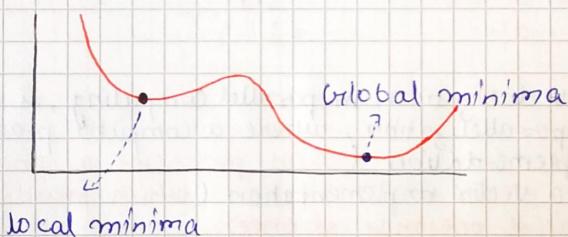
Let's say we have a cost function  $J(w)$  which we want to minimize. What gradient descent does is, it repeatedly performs the following updating steps:



until the algorithm converges at Global optimum.

★ **Global optima** → A global minimum of a function is a point where the function value is smaller than all other feasible points.

★ **Local optima** → A local minimum of a function is a point where the function value is smaller than at nearby points, but possibly greater than at a distant point.



→ why this update rule makes sense is that :

- ① If we are on the rising side (+ve slope) of the function, we would have a +ve slope, therefore  $w = w - (\dots)$  i.e. the value will be reduced.
- ② If we are on the falling side (-ve slope) of the function, we would have a -ve slope, therefore  $w = w + (\dots)$  i.e. the value will be increased.

⇒ The final update rule becomes

$$\boxed{w := w - \alpha \frac{dJ(w, b)}{dw} \rightarrow \frac{\partial J(w, b)}{\partial w}}$$

$$b := b - \alpha \frac{dJ(w, b)}{db} \rightarrow \frac{\partial J(w, b)}{\partial b}$$

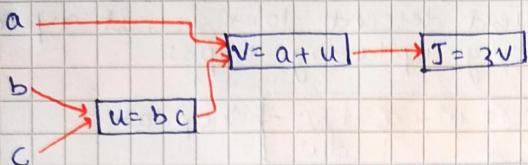
In case of multivariable, we use  $\nabla(\alpha)$  instead of  $d(\alpha)$

→ Computation Graph → It is defined as a directed graph where the nodes corresponds to mathematical operations. Computational graphs are a way of expressing and evaluating a mathematical expression.

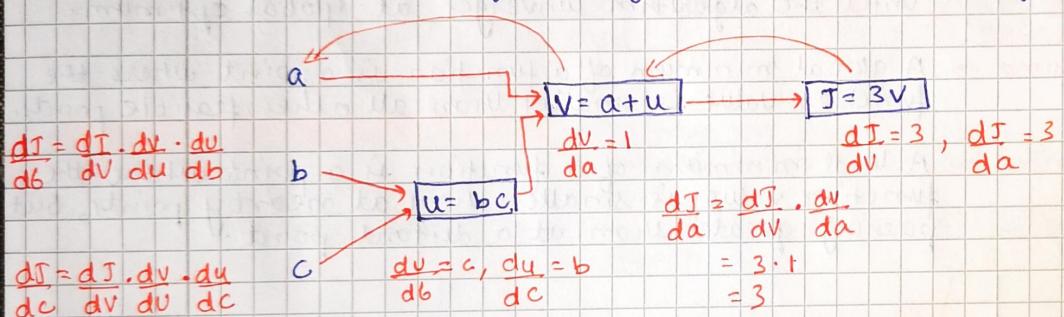
Eg.  $J(a, b, c) = 3(a + bc)$

$$\begin{array}{c} u \\ \overbrace{\quad\quad}^u \\ v \\ \overbrace{\quad\quad}^v \\ J \end{array}$$

$$\begin{aligned} u &= bc \\ v &= a + u \\ J &= 3v \end{aligned}$$



Backward Pass (Back propagation) → In the backward pass, our intention is to calculate the gradient for each input with respect to the final output. These gradients are essential for training the neural network using gradient descent.



→ Vectorization → Automatic vectorization, in parallel computing, is a special case of automatic parallelization, where a computer program is converted from a scalar implementation (which processes a single pair of operands at a time) to a vector implementation (which processes one operation on multiple pairs of operands at once).

\* whenever possible avoid "for-loops". Instead use in-built functions.

→ Calculating the derivatives of logistic regression (with and without Vectorization)

$$\begin{aligned} J &= 0, dw_1 = 0, dw_2 = 0, db = 0 \\ \text{for } i &= 1 \text{ to } m: \\ z^i &= w^T x^i + b \\ a^i &= \sigma(z^i) \\ J &+= -[y^i \log(\hat{y}^i) + (1-y^i) \log(1-\hat{y}^i)] \\ dz^i &= a^i(1-a^i) \\ dw_1 &+= x_1^i dz^i \\ dw_2 &+= x_2^i dz^i \\ db &+= dz^i \end{aligned}$$

$$\begin{aligned} J &= J/m, dw_1 = dw_1/m, dw_2 = dw_2/m \\ db &= db/m \end{aligned}$$

$$\begin{aligned} J &= 0, dw = np.zeros((m, 1)), db = 0 \\ \text{for } i &= 1 \text{ to } m: \\ z^i &= w^T x^i + b \\ a^i &= \sigma(z^i) \\ J &+= \dots \\ dz^i &= a^i(1-a^i) \\ dw &+= x^i dz^i \\ db &+= dz^i \end{aligned}$$

$$J = J/m, dw_1 = m, db = db/m$$

➤ **Broadcasting** → The term broadcasting refers to how numpy treats arrays with different dimensions during arithmetic operations which leads to certain constraints, the smaller array is broadcast across the larger array so that they have compatible shapes.